# Community Detection: Overlapping Communities

CS224W: Social and Information Network Analysis
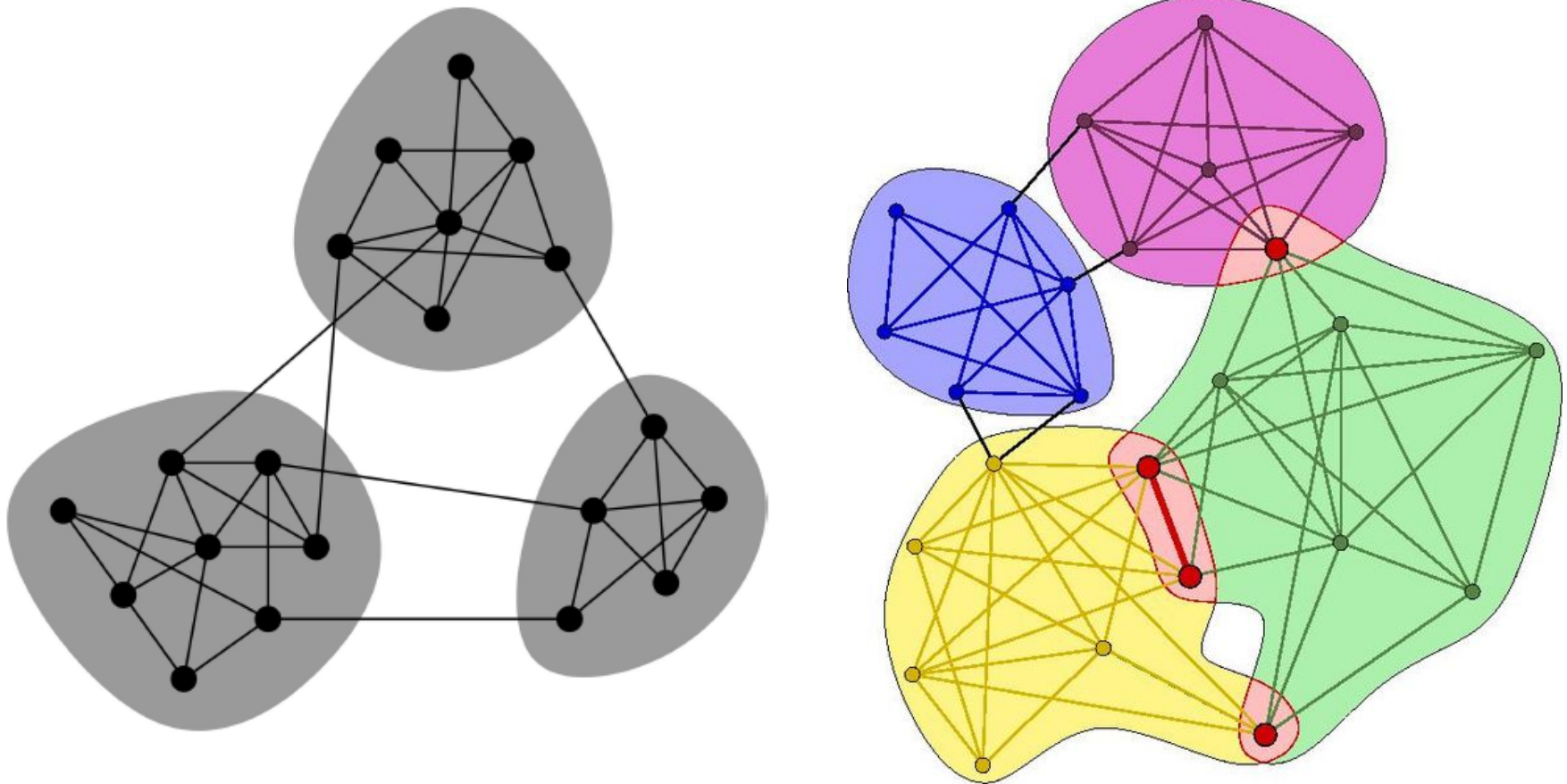Jure Leskovec, Stanford University
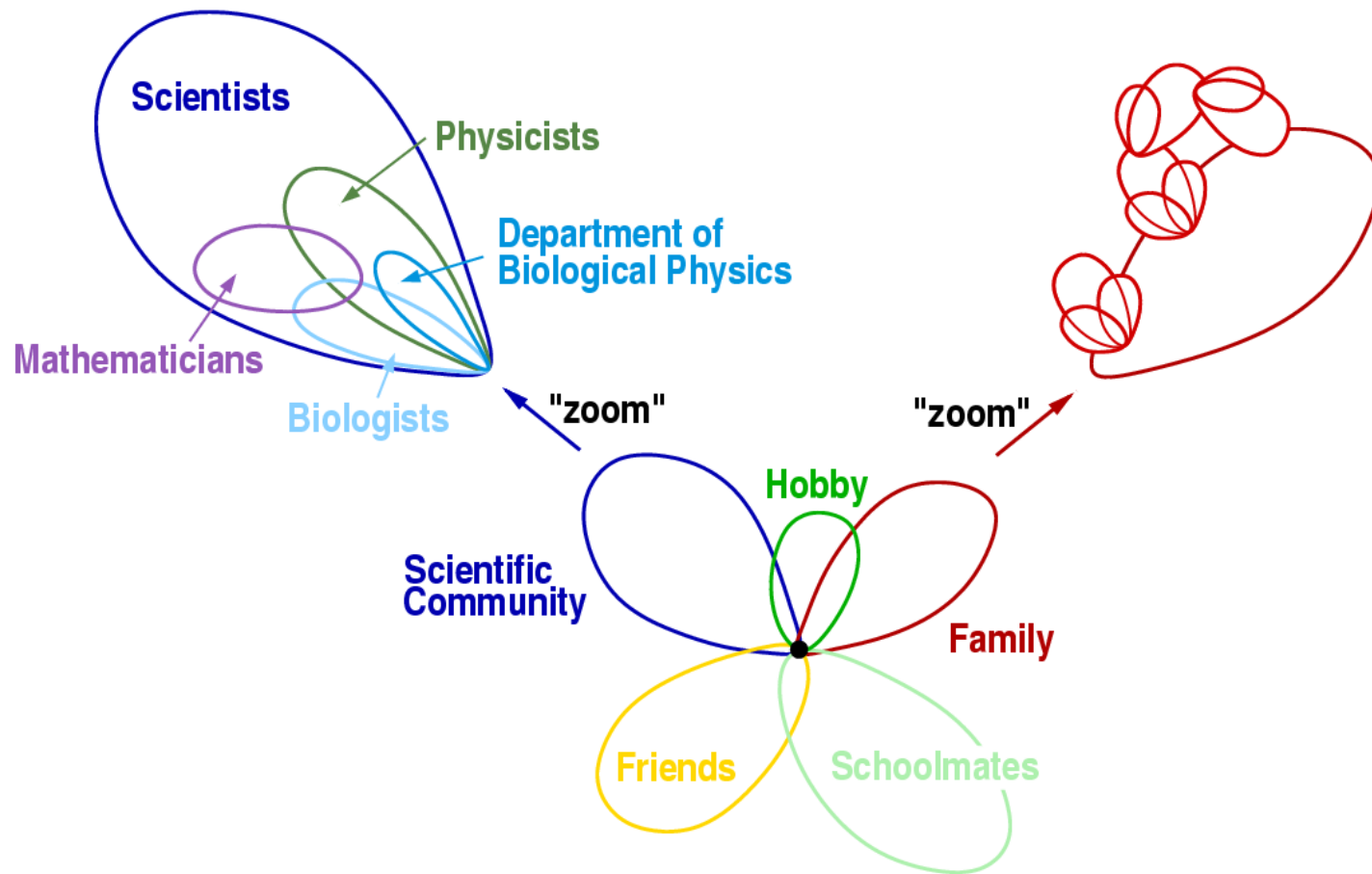http://cs224w.stanford.edu
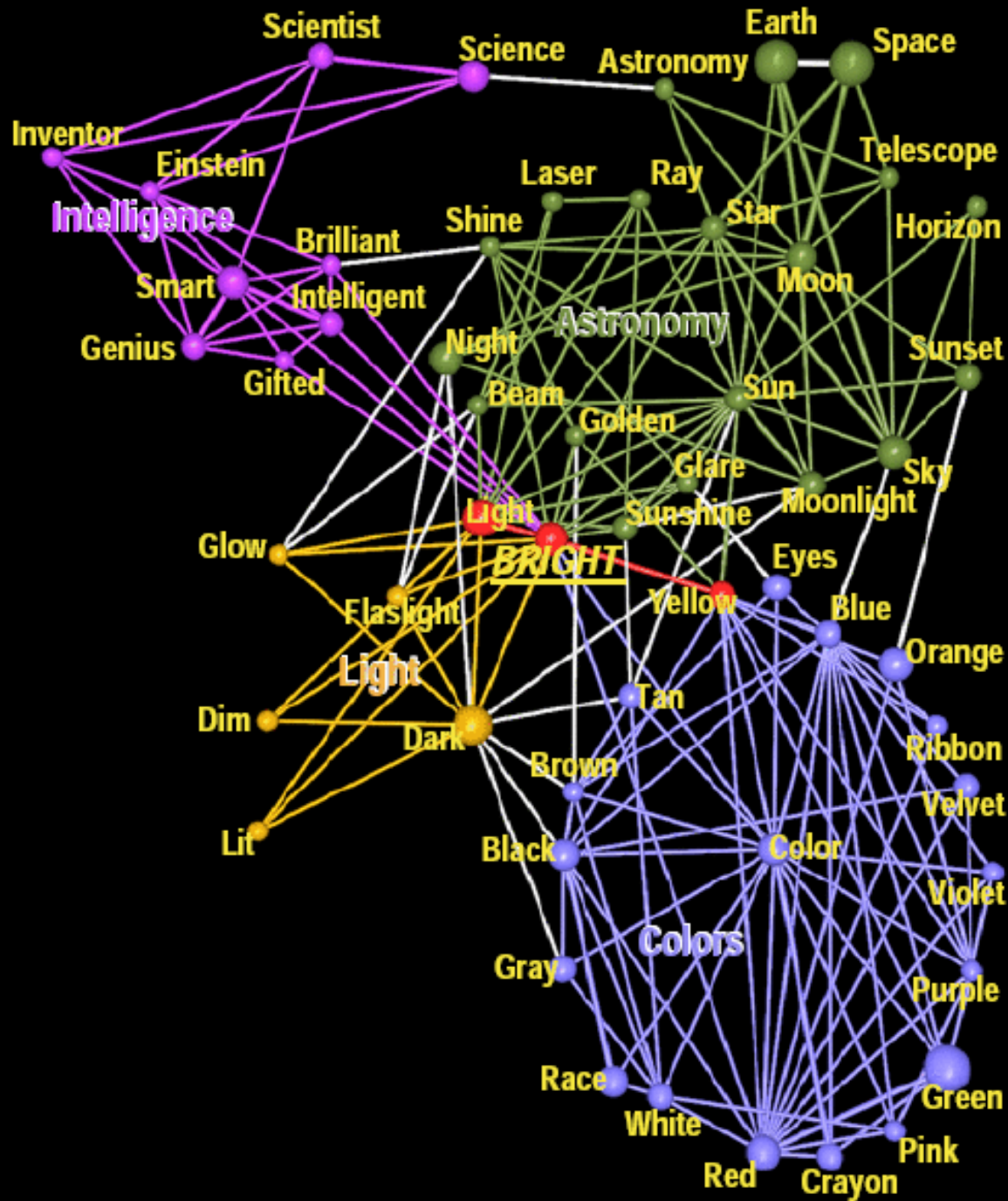
# Overlapping Communities

- **Non-overlapping vs. overlapping  communities**

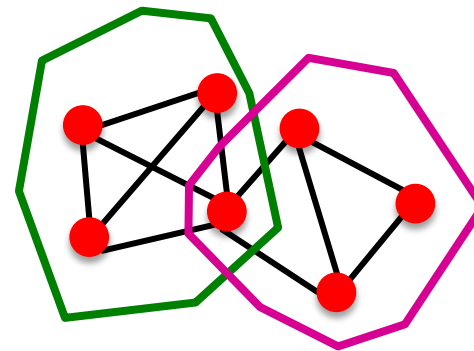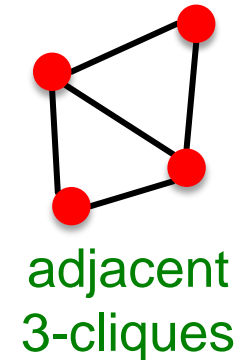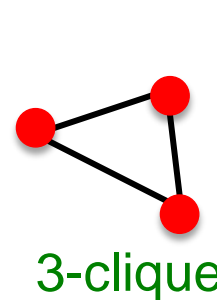# Overlaps of Social Circles

- **A node belongs to many social circles**

# Clique Percolation Method (CPM)

- **Two nodes belong to the same community if they can be connected through adjacent $k$-cliques:**

    - **$k$-clique:**
        - Fully connected graph on $k$ nodes

    - **Adjacent $k$-cliques:**
        - overlap in $k$-$1$ nodes

- **$k$-clique community**
    - Set of nodes that can be reached through a sequence of adjacent $k$-cliques
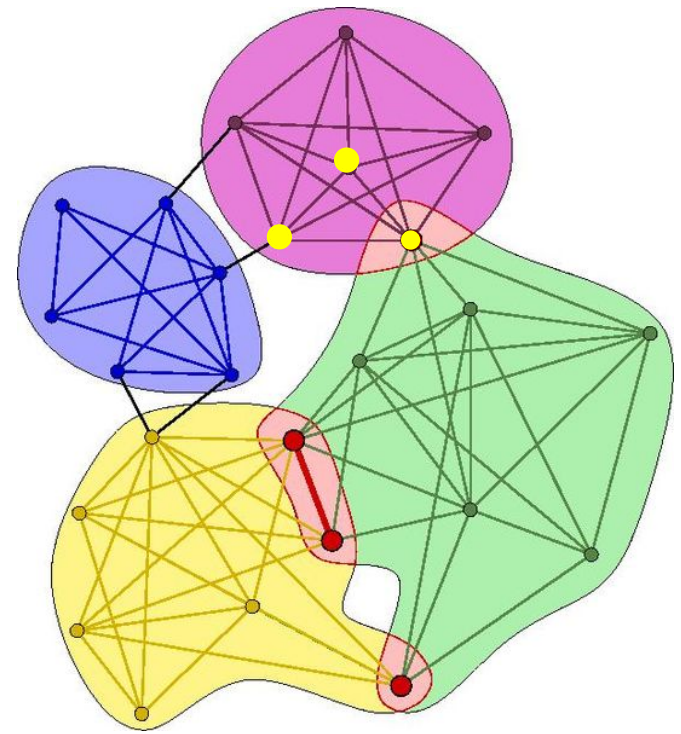
3-clique

adjacent
3-cliques

# Clique Percolation Method (CPM)

- **Two nodes belong to the same community if they can be connected through adjacent $k$-cliques:**
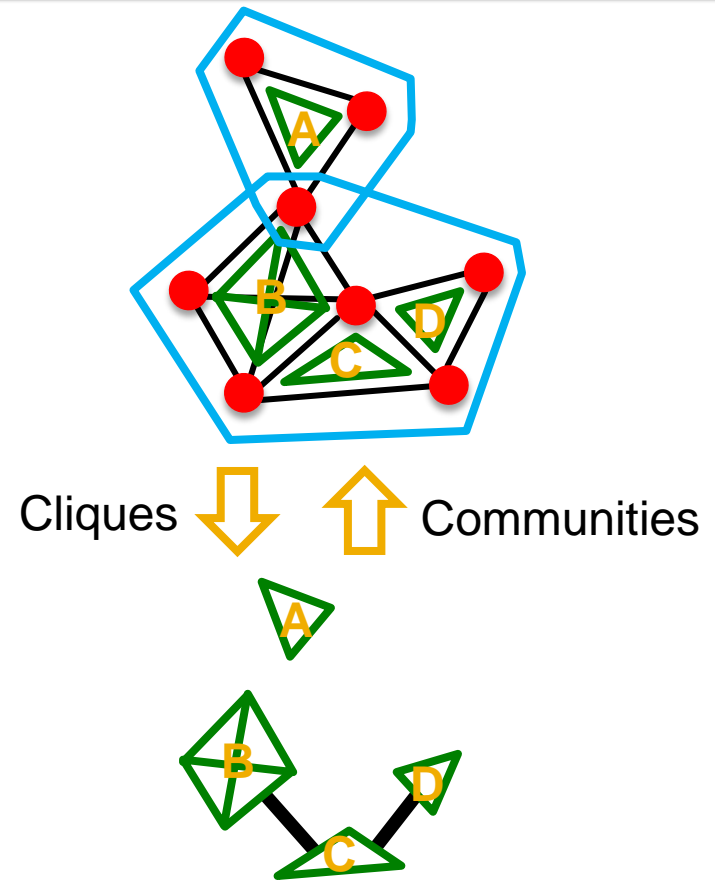
4-clique

# CPM: Steps

- **Clique Percolation Method:**
  - **Find maximal-cliques** (not $k$-cliques!)
  - **Clique overlap graph:**
    - Each clique is a node
    - Connect two cliques if they overlap in at least $k$-$1$ nodes
  - **Communities:**
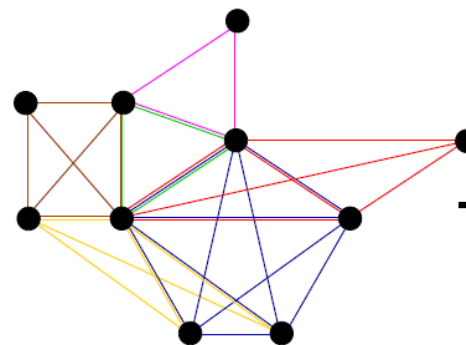    - Connected components of the clique overlap matrix
- **How to set $k$?**
  - Set $k$ so that we get the "richest" (most widely distributed cluster sizes) community structure

Cliques ⇩   ⇧ Communities

# CPM method: Example

- Start with graph
- Find maximal cliques
- Create clique overlap matrix
- Threshold the matrix at value *k-1*
  - If $a_{ij} < k\text{-}1$ set 0
- Communities are the connected components of the thresholded matrix



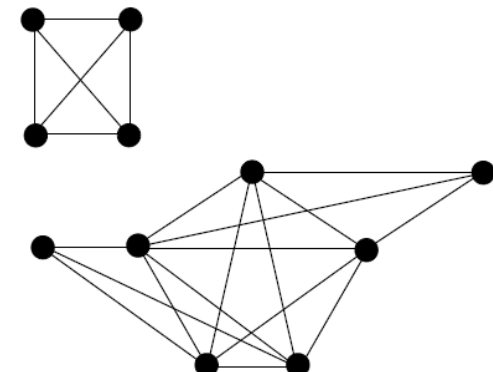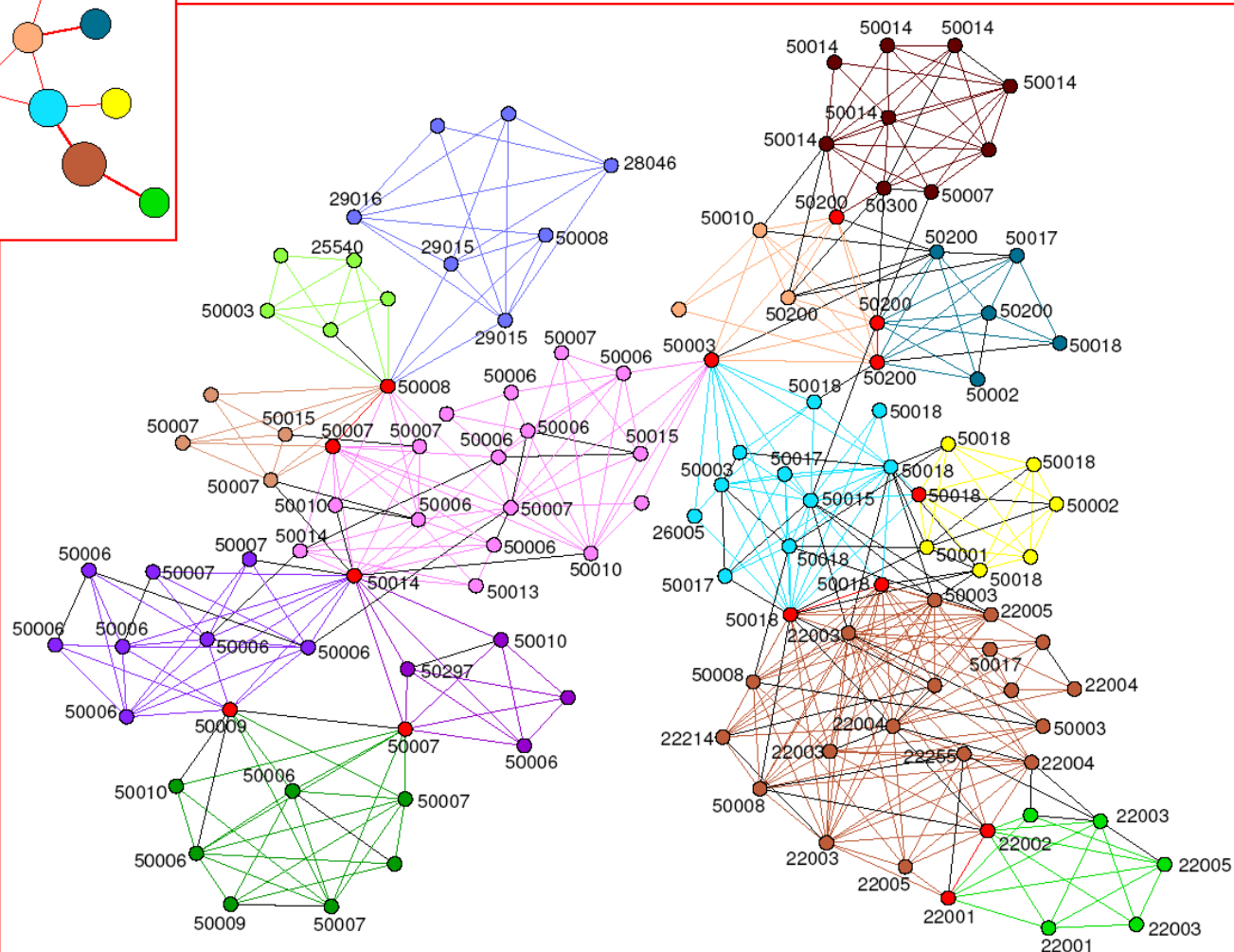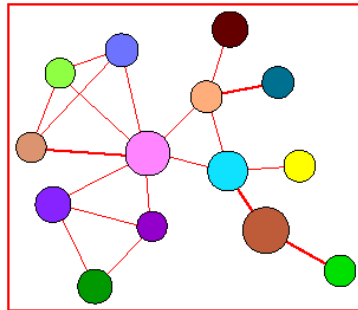**(1) Graph**

**(2) Clique overlap matrix**

k=4

**(3) Thresholded matrix at 3**
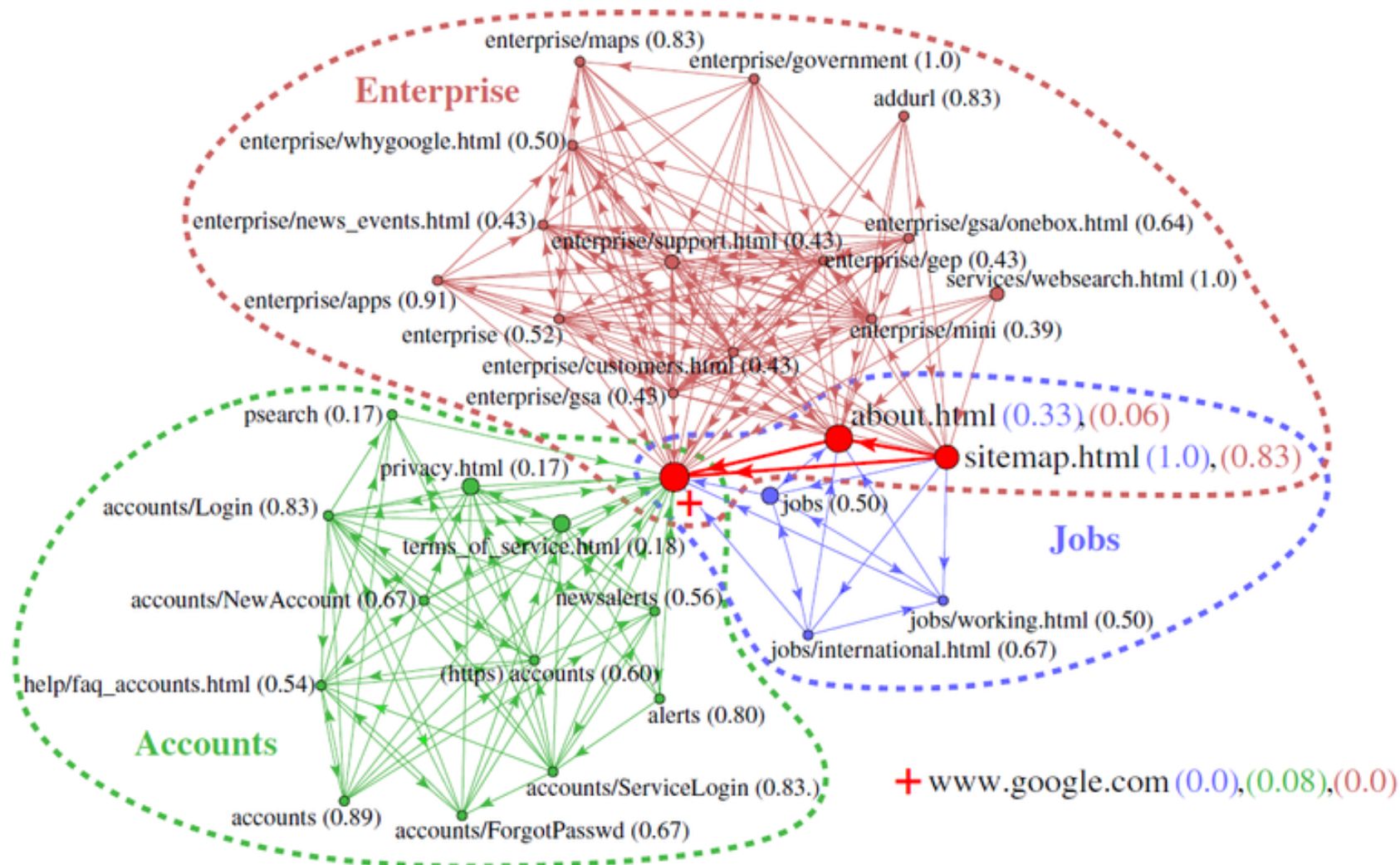
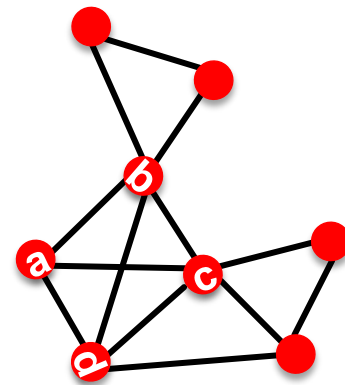**(4) Communities (connected components)**

# Example: Phone-Call Network



Communities in a "tiny" part of a phone call network of 4 million users
[Palla et al., '07]

# Example: Website



Enterprise
enterprise/maps (0.83)
enterprise/government (1.0)
addurl (0.83)
enterprise/whygoogle.html (0.50)
enterprise/news_events.html (0.43)
enterprise/gsa/onebox.html (0.64)
enterprise/support.html (0.43)
enterprise/gep (0.43)
services/websearch.html (1.0)
enterprise/apps (0.91)
enterprise (0.52)
enterprise/mini (0.39)
enterprise/customers.html (0.43)
enterprise/gsa (0.43)
about.html (0.33),(0.06)
psearch (0.17)
sitemap.html (1.0),(0.83)
privacy.html (0.17)
accounts/Login (0.83)
jobs (0.50)
Jobs
terms_of_service.html (0.18)
accounts/NewAccount (0.67)
newsalerts (0.56)
jobs/working.html (0.50)
jobs/international.html (0.67)
(https) accounts (0.60)
help/faq_accounts.html (0.54)
alerts (0.80)
Accounts
accounts/ServiceLogin (0.83.)
www.google.com (0.0),(0.08),(0.0)
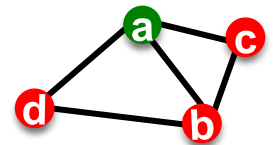accounts (0.89)
accounts/ForgotPasswd (0.67)

# How to Find Maximal Cliques?

- **No nice way, NP-hard combinatorial problem**
- **Maximal clique:** clique that can't be extended
  - {a,b,c} is a clique but not maximal clique
  - {a,b,c,d} is maximal clique
- **Algorithm:** Sketch
  - Start with a seed node
  - Expand the clique around the seed
  - Once the clique cannot be further expanded we found the maximal clique
  - **Note:**
    - This will generate the same clique multiple times

# How to Find Maximal Cliques?

- Start with a seed vertex "a"
- **Goal:** Find the maximal clique Q "a" belongs to
  - **Observation:**
    - If some "x" belongs to Q then it is a member of "a"
      - **Why?** If $a,x \in Q$ but not a–x, then Q is not a clique!
- **Recursive algorithm:**
  - Q … current clique
  - R … candidate vertices to expand the clique to
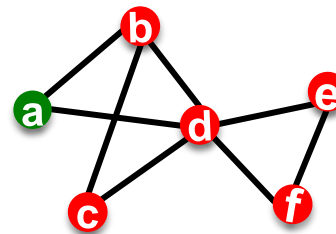- **Example:** Start with "a" and expand around it

| Q= | {a} | {a,b} | {a,b,c} | bktrack | {a,b,d} |
|---|---|---|---|---|---|
| R= | {<u>b</u>,c,d} | {b,c,d} | {d}$\cap\Gamma$(c)={} | | {c}$\cap\Gamma$(d)={} |
| | | $\cap\Gamma$(b)={<u>c</u>,d} | | | |

Steps of the recursive algorithm

$\Gamma$(u)…neighbor set of u

# How to Find Maximal Cliques?

- Q … current clique
- R … candidate vertices

- **Expand(R,Q)**

  - **while** R $\neq$ { }

    - p = vertex in R
    - $Q_p$ = Q $\cup$ {p}
    - $R_p$ = R $\cap$ $\Gamma$(p)
    - **if** $R_p$ $\neq$ { }: Expand($R_p,Q_p$)
      **else:** output $Q_p$
    - R = R - {p}

**Start: Expand(V, {})**
R={a,…f}, Q={}
p = {a}
$Q_p$ = {a}
$R_p$ = {b,d}
**Expand($R_p$, Q):**
  R = {b,d}, Q={a}
  p = {b}
  $Q_p$ = {a,b}
  $R_p$ = {d}
  **Expand($R_p$, Q):**
    R = {d}, Q={a,b}
    p = {d}
    $Q_p$ = {a,b,d}
    $R_p$ = {} : **output {a,b,d}**
  p = {d}
  $Q_p$ = {a,d}
  $R_p$ = {b}
  **Expand($R_p$, Q):**
    R = {b}, Q={a,d}
    p = {b}
    $Q_p$ = {a,d}
    $R_p$ = {} : **output {a,d,b}**

- Q ... current clique
- R ... candidate vertices
- **Expand(R,Q)**
  - **while** $R \neq \{\}$
    - $p$ = vertex in R
    - $Q_p = Q \cup \{p\}$
    - $R_p = R \cap \Gamma(p)$
    - **if** $R_p \neq \{\}$: $\text{Expand}(R_p, Q_p)$
      **else:** output $Q_p$
    - $R = R - \{p\}$

**Start: Expand(V, {})**
 $R=\{a,\ldots f\}$, $Q=\{\}$
 $p = \{b\}$
 $Q_p = \{b\}$
 $R_p = \{a,c,d\}$
**Expand($R_p$, Q):**
  $R = \{a,c,d\}$, $Q=\{b\}$
  $p = \{a\}$
  $Q_p = \{b,a\}$
  $R_p = \{d\}$
  **Expand($R_p$, Q):**
   $R = \{d\}$, $Q=\{b,a\}$
   $p = \{d\}$
   $Q_p = \{b,a,d\}$
   $R_p = \{\}$ : **output {b,a,d}**
 $p = \{c\}$
 $Q_p = \{b,c\}$
 $R_p = \{d\}$
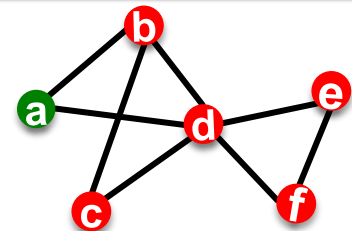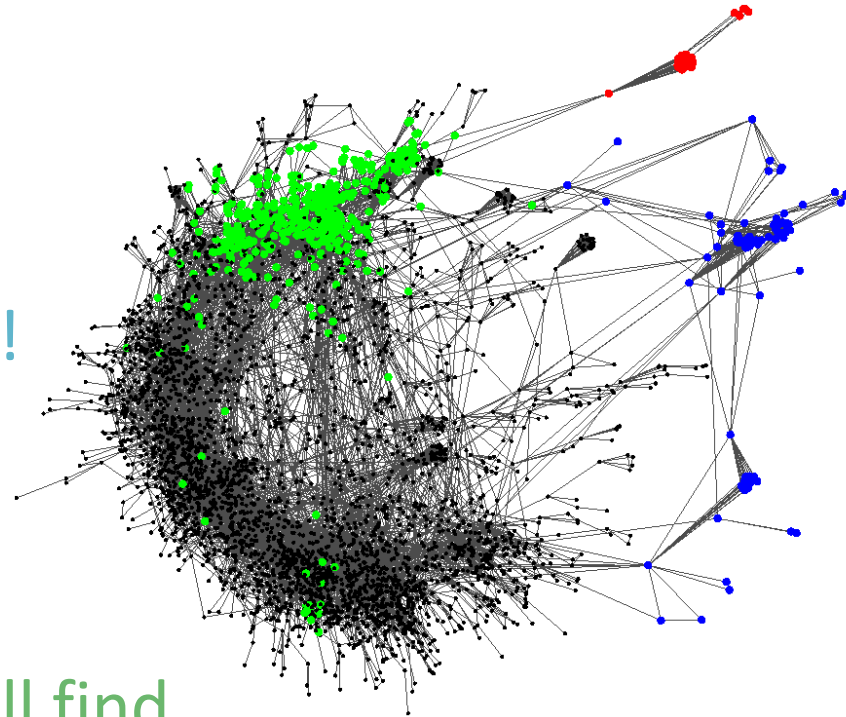 **Expand($R_p$, Q):**
   $R = \{d\}$, $Q=\{b,c\}$
   $p = \{d\}$
   $Q_p = \{b,c,d\}$
   $R_p = \{\}$ : **output {b,c,d}**

# How to Find Maximal Cliques?

- **How to prevent maximal cliques to be generated multiple times?**

  - Only output cliques that are lexicographically minimum
    - $\{a,b,c\} < \{b,a,c\}$

  - **Even better:** Only expand to the nodes higher in the lexicographical order

**Start: Expand(V, {})**
$R = \{a,\ldots f\}$, $Q = \{\}$
$p = \{a\}$
$Q_p = \{a\}$
$R_p = \{b,d\}$
**Expand($R_p$, Q):**
  $R = \{b,d\}$, $Q = \{a\}$
  $p = \{b\}$
  $Q_p = \{a,b\}$
  $R_p = \{d\}$
  **Expand($R_p$, Q):**
    $R = \{d\}$, $Q = \{a,b\}$
    $p = \{d\}$
    $Q_p = \{a,b,d\}$
    $R_p = \{\}$ : **output {a,b,d}**
  $p = \{d\}$
  $Q_p = \{a,d\}$   **Don't expand**
  $R_p = \{b\}$     **b < d**

# How to Model Networks with Communities?

# Reflections: Finding Communities

- **Let's rethink what we are doing...**
  - Given a network
  - Want to find communities!
- **Need to:**
  - Formalize the notion of a community
  - Need an algorithm that will find sets of nodes that are "good" communities
- **More generally:**
  - **How to think about clusters in large networks?**

# Clustering Objective Functions
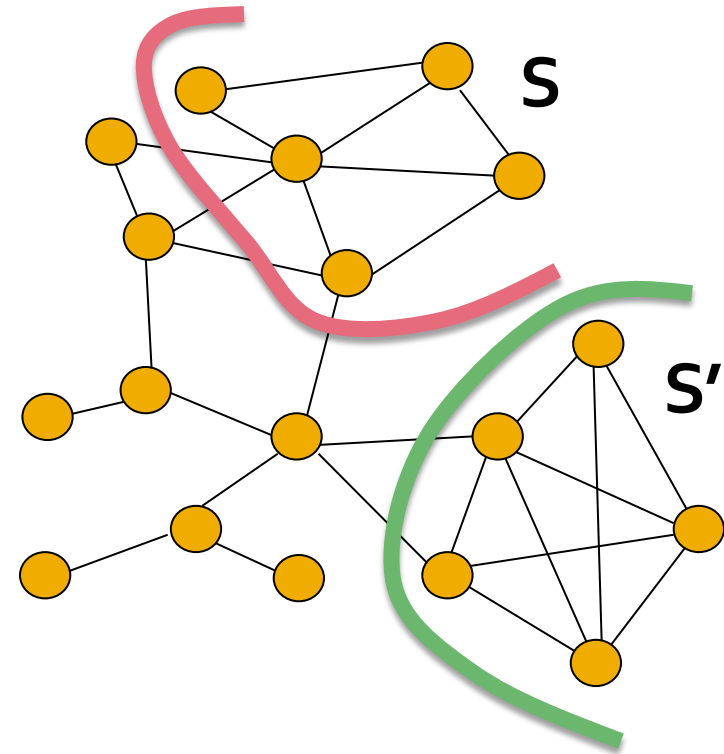
**What is a good cluster?**
- Many edges internally
- Few pointing outside

Formally, **conductance:**

$$\phi(S) = \frac{\sum_{i \in S, j \notin S} A_{ij}}{\min\{A(S), A(\overline{S})\}}$$

Where: A(S)....volume $\quad A(S) = \sum_{i \in S} \sum_{j \in V} A_{ij}$

**Small Φ(S) corresponds to good clusters**

# Community Score

- **How community like is a set of nodes?**
- A good cluster *S* has
  - Many edges internally
  - Few edges pointing outside
- Simplest objective function: **Conductance**

$$\phi(S) = \frac{|\{(i,j) \in E; i \in S, j \notin S\}|}{\sum_{s \in S} d_s}$$

**Small** **conductance** corresponds to good clusters

# Network Community Profile Plot

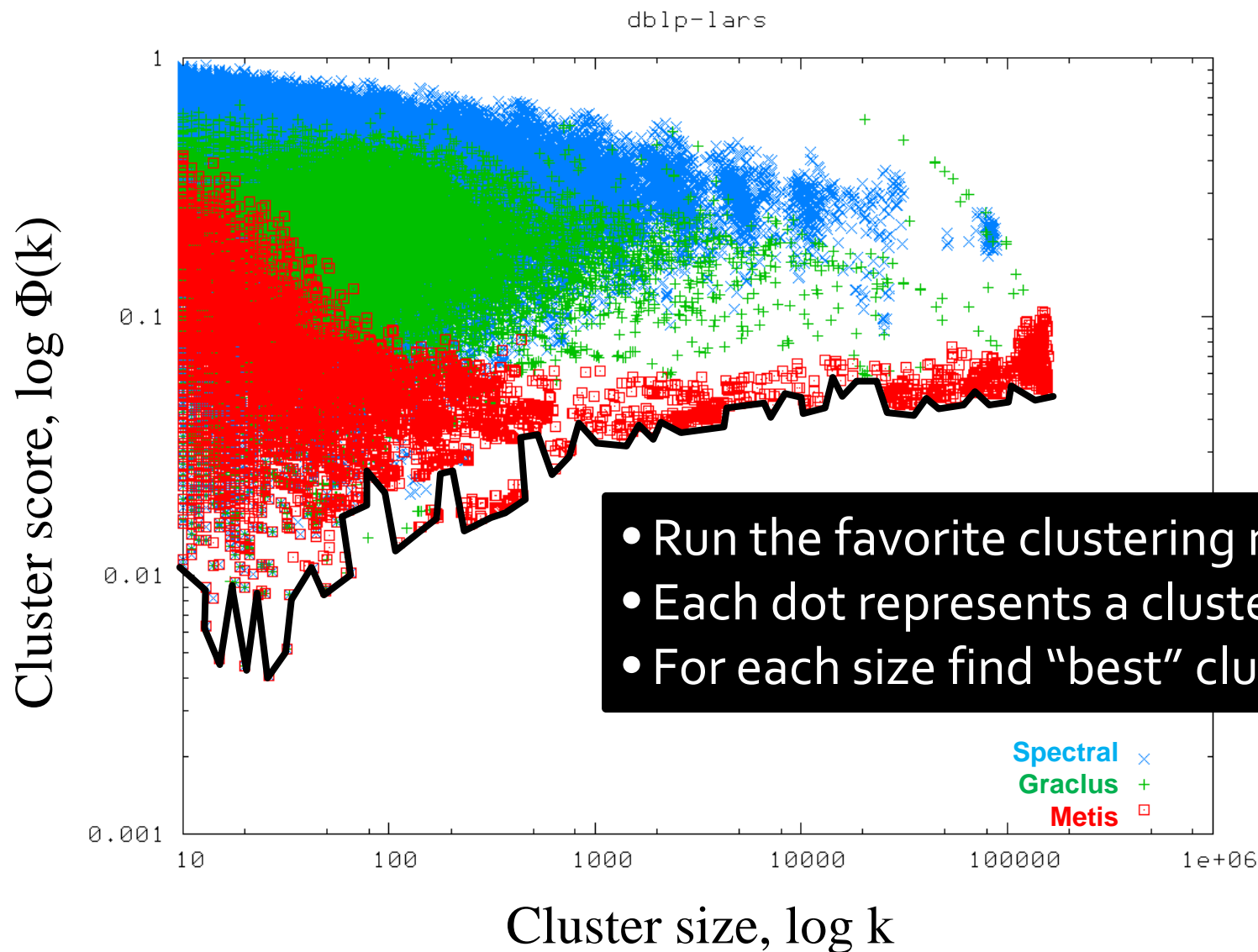- Define:

Network community profile (**NCP**) plot

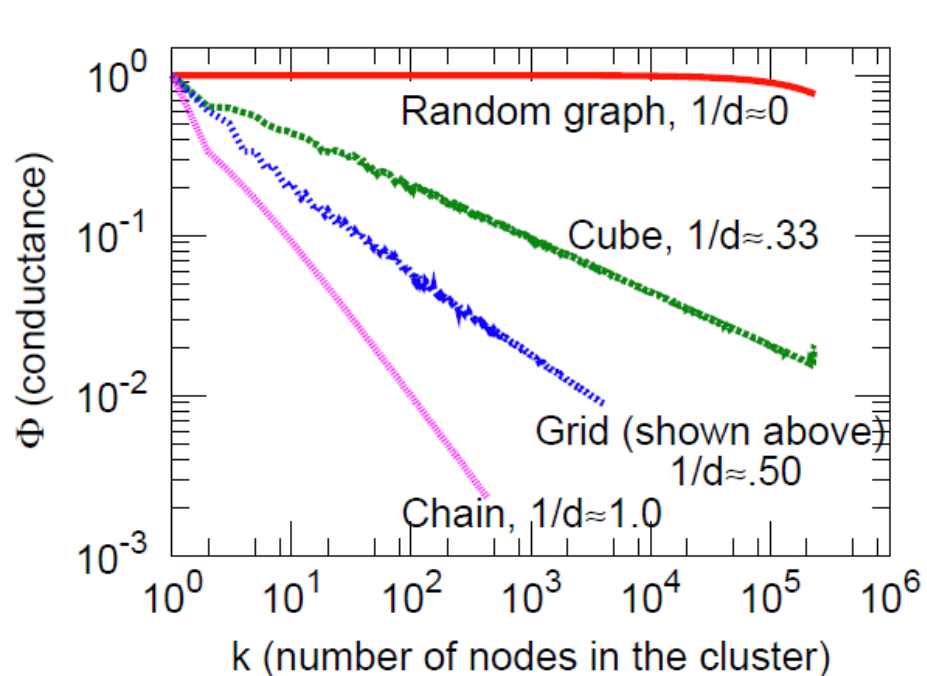Plot the score of **best** community of size *k*

$$\Phi(k) = \min_{S \subset V, |S|=k} \phi(S)$$

dblp-lars

Cluster score, log $\Phi(k)$

Cluster size, log k

- Run the favorite clustering method
- Each dot represents a cluster
- For each size find "best" cluster

**Spectral** ×
**Graclus** +
**Metis** □

# NCP Plot: Meshes

- ## **Meshes, grids, dense random graphs:**



d-dimensional meshes



California road network

# NCP plot: Network Science
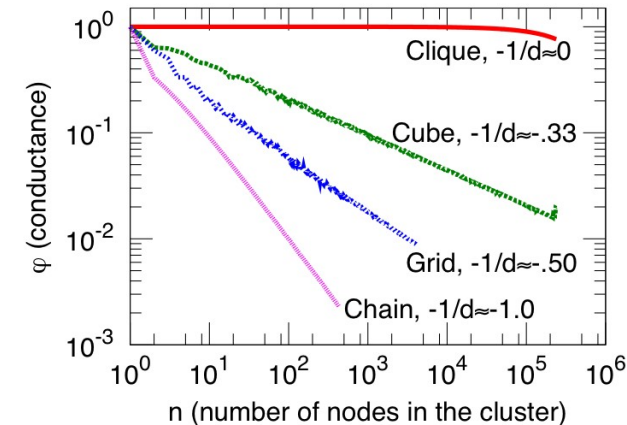
- ## Collaborations between scientists in networks
  [Newman, 2005]

# Natural Hypothesis

## Natural hypothesis about NCP:

- NCP of real networks slopes downward
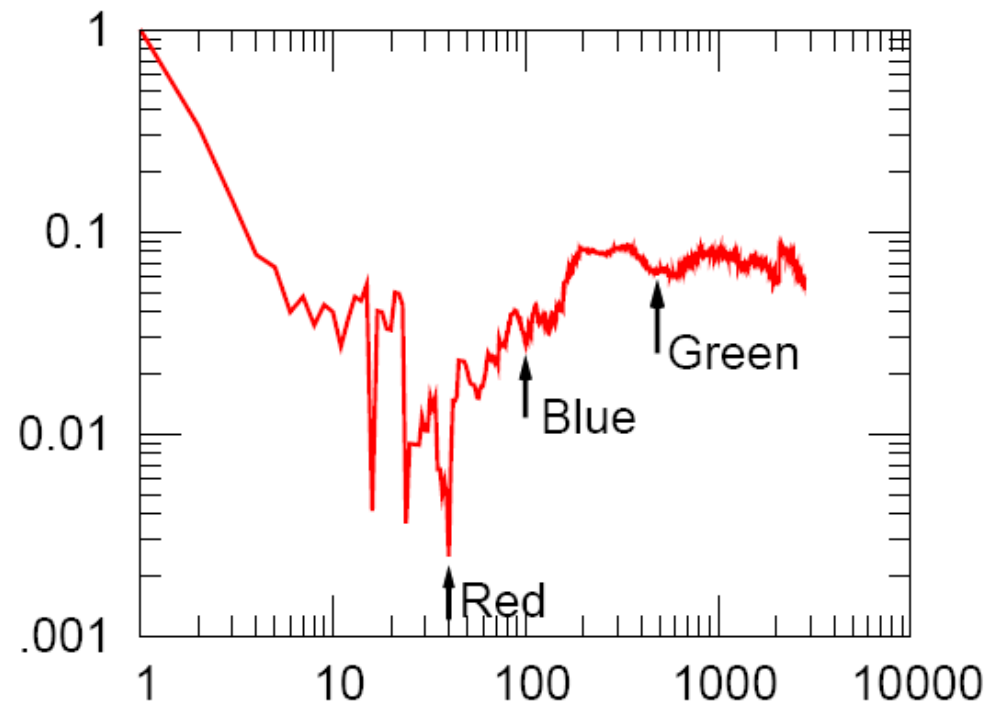- Slope of the NCP corresponds to the "dimensionality" of the network



Clique, -1/d≈0
Cube, -1/d≈-.33
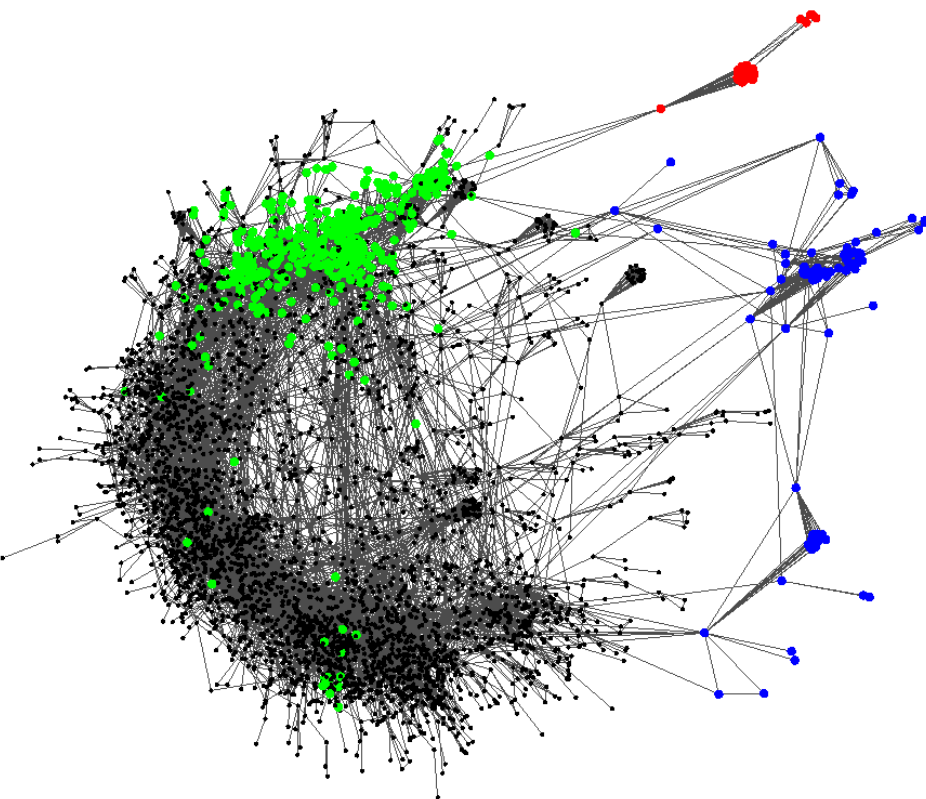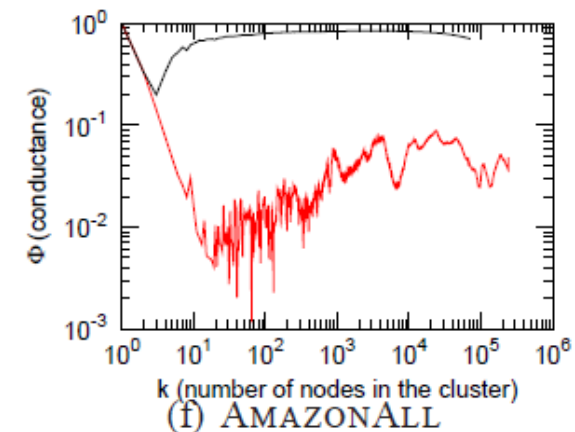Grid, -1/d≈-.50
Chain, -1/d≈-1.0

φ (conductance) vs n (number of nodes in the cluster)

### What about large networks?

| • Social nets | Nodes | Edges | Description |
|---|---|---|---|
| LiveJournal | 4,843,953 | 42,845,684 | Blog friendships [5] |
| Epinions | 75,877 | 405,739 | Trust network [28] |
| CA-DBLP | 317,080 | 1,049,866 | Co-authorship [5] |
| • Information (citation) networks | | | |
| Cit-hep-th | 27,400 | 352,021 | Arxiv hep-th [14] |
| AmazonProd | 524,371 | 1,491,793 | Amazon products [8] |
| • Web graphs | | | |
| Web-Google | 855,802 | 4,291,352 | Google web graph |
| Web-WT10G | 1,458,316 | 6,225,033 | TREC WT10G |
| • Bipartite affiliation (authors-to-papers) networks | | | |
| Atp-DBLP | 615,678 | 944,456 | DBLP [21] |
| Atm-Imdb | 2,076,978 | 5,847,693 | Actors-to-movies |
| • Internet networks | | | |
| AsSkitter | 1,719,037 | 12,814,089 | Autonom. sys. |
| Gnutella | 62,561 | 147,878 | P2P network [29] |

# Large Networks: Very Different

**Typical example:** General Relativity collaborations (n=4,158, m=13,422)
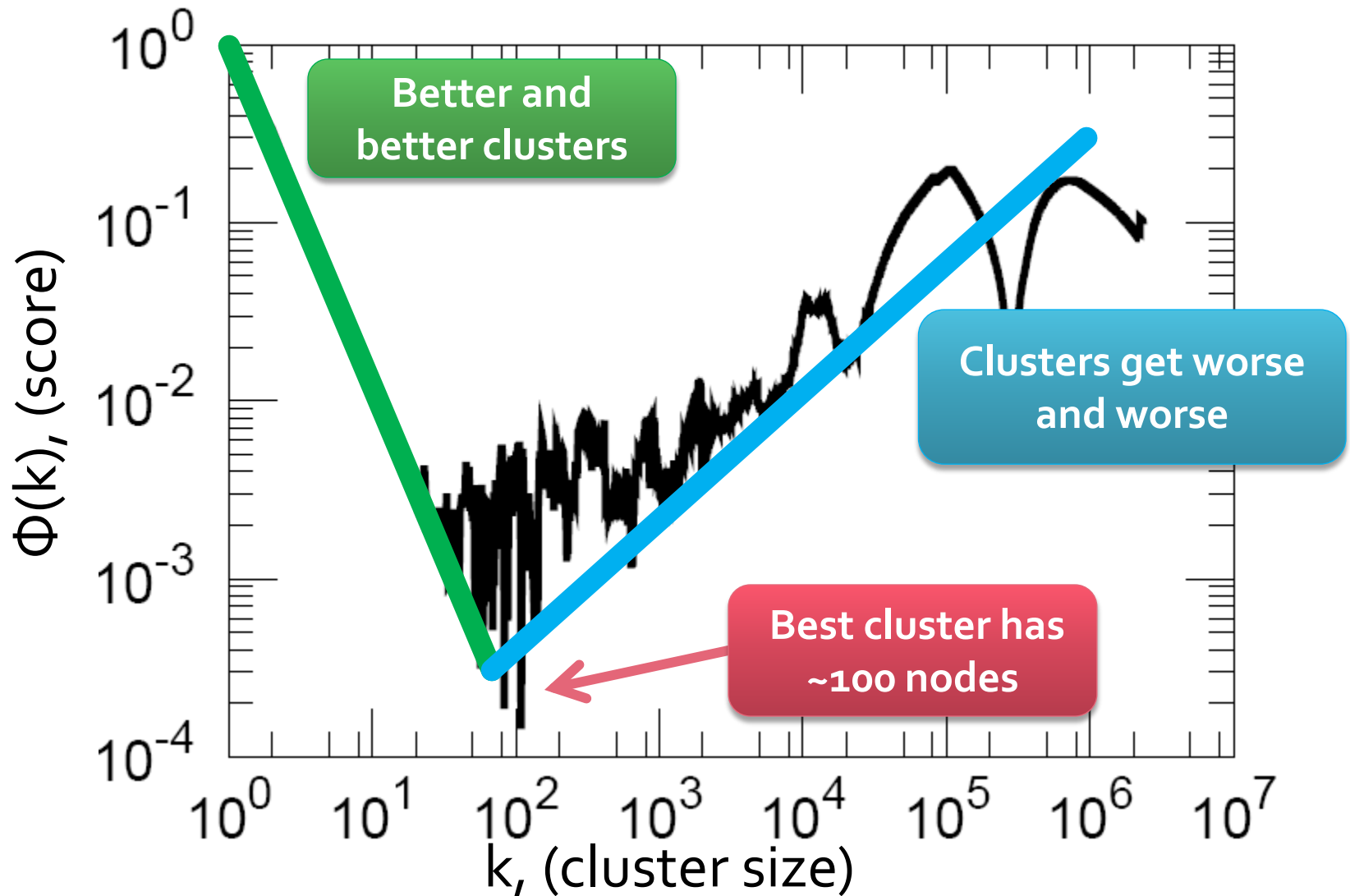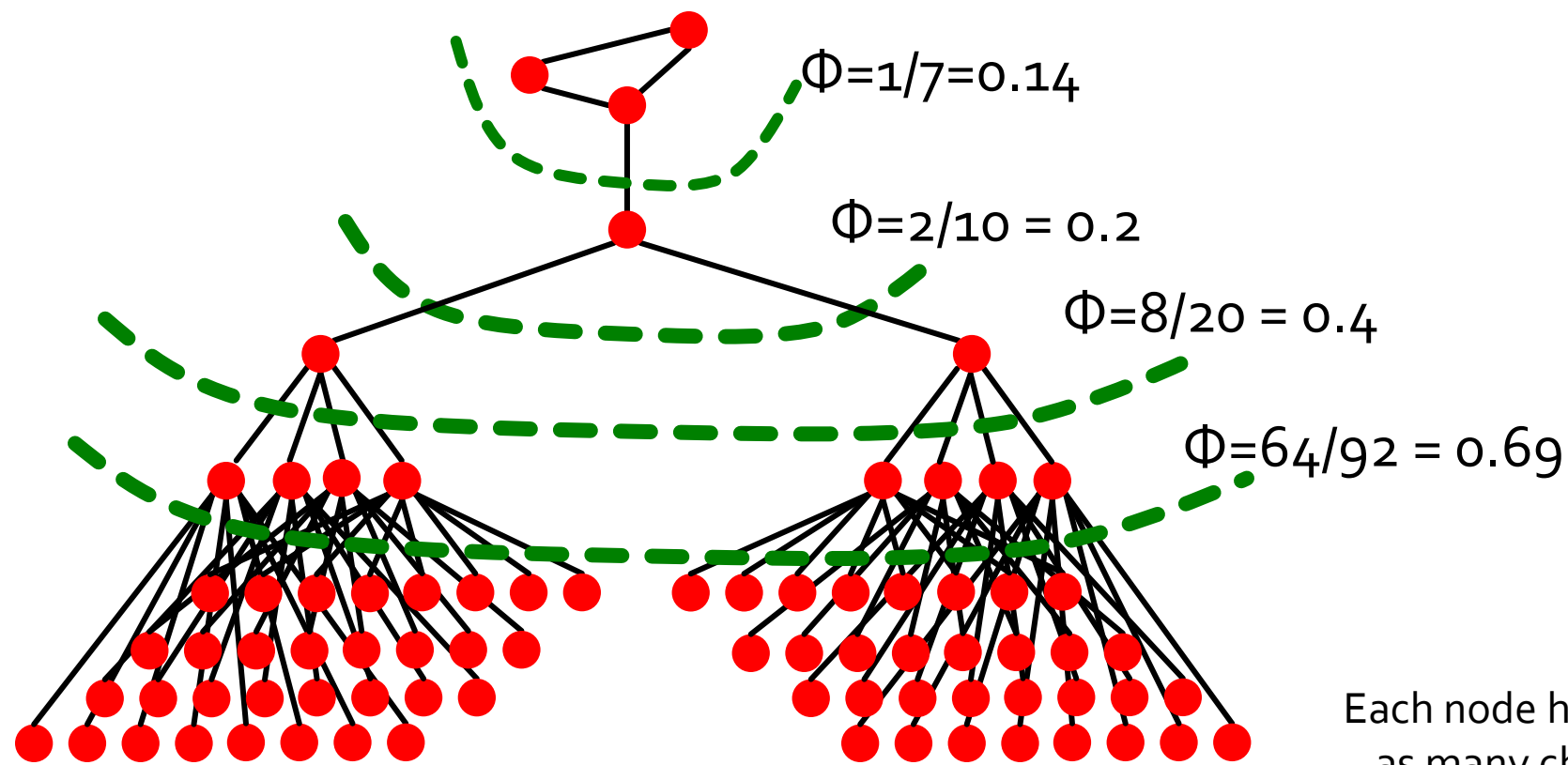
# More NCP Plots of Networks



(a) LiveJournal01

(b) Messenger-DE

(c) AtP-DBLP

(d) Cit-hep-th

(e) Web-Google

(f) AmazonAll

- As clusters grow the number of edges inside grows **slower** that the number crossing



$\Phi=1/7=0.14$

$\Phi=2/10 = 0.2$

$\Phi=8/20 = 0.4$

$\Phi=64/92 = 0.69$
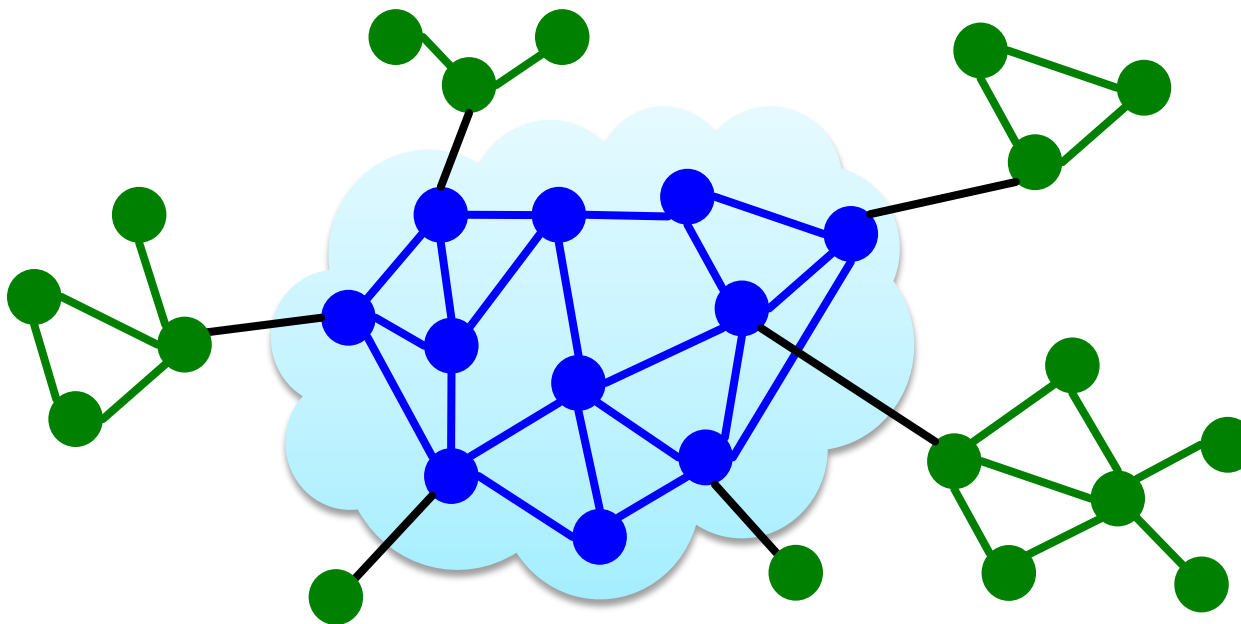
Each node has twice as many children

28

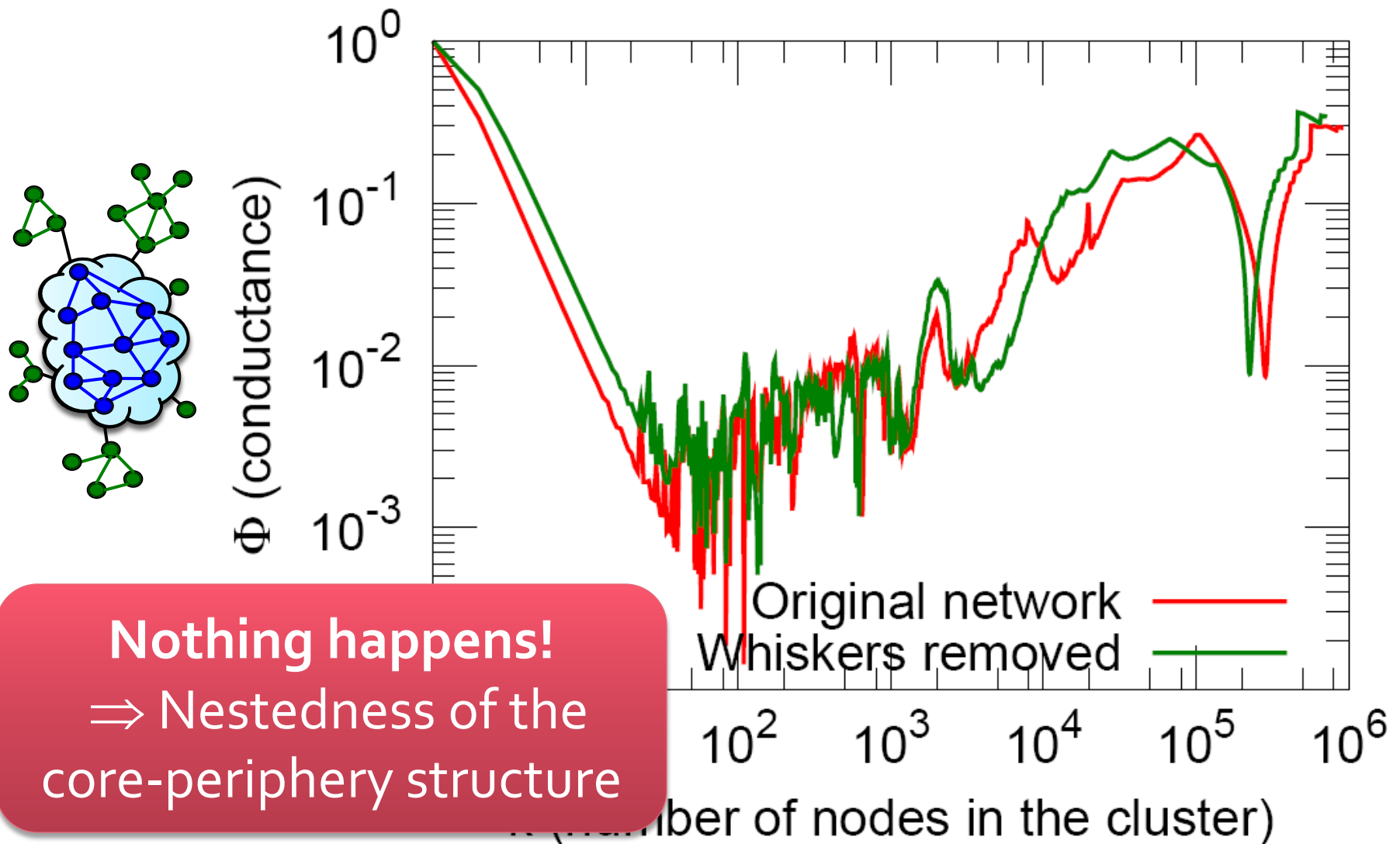- Empirically we note that best clusters are barely connected to the network



NCP plot

⇒ **Core-periphery structure**

**Nothing happens!**
$\Rightarrow$ Nestedness of the core-periphery structure

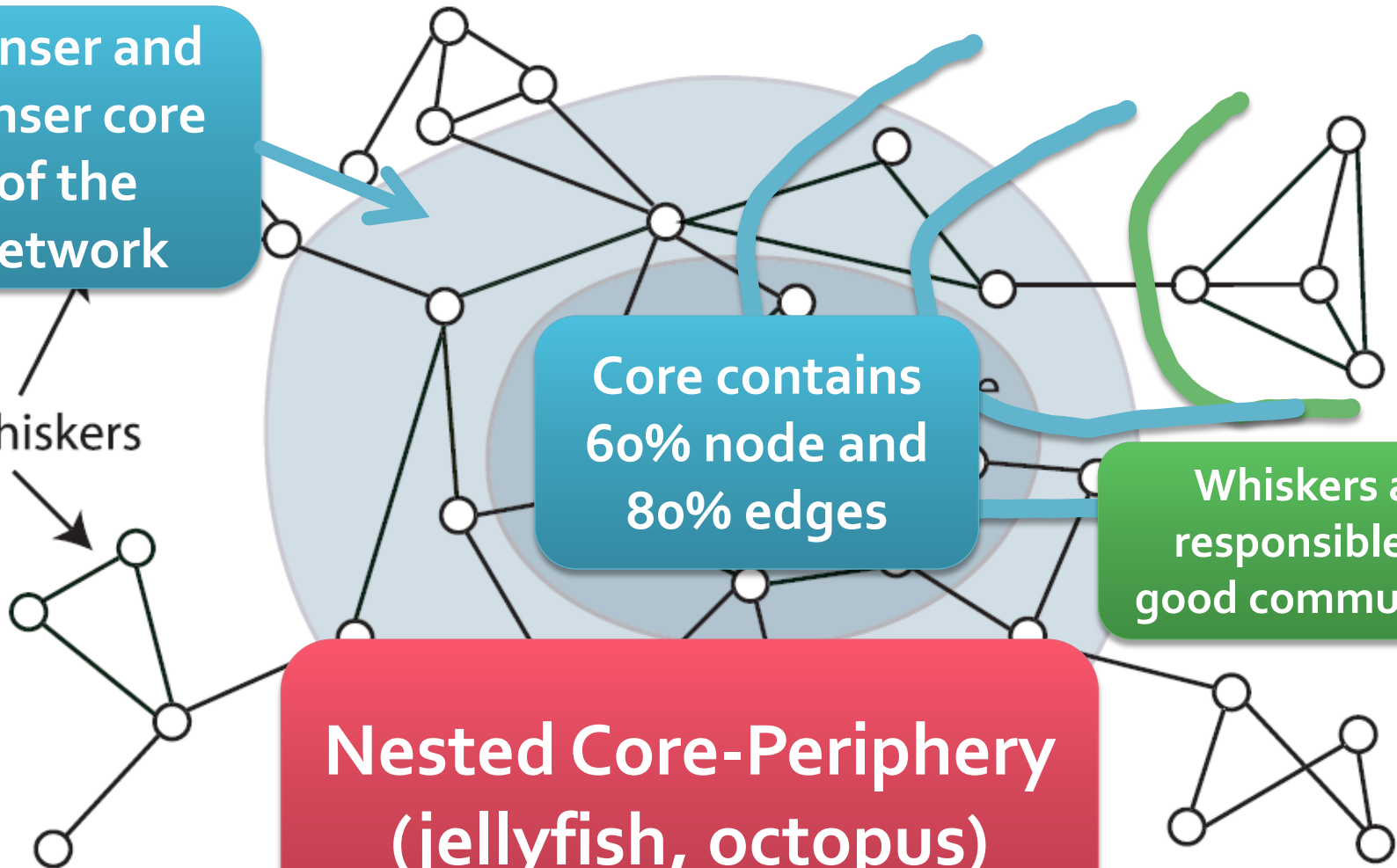# Suggested Network Structure

**Denser and denser core of the network**

Whiskers

**Core contains 60% node and 80% edges**

**Whiskers are responsible for good communities**

**Nested Core-Periphery (jellyfish, octopus)**

# Communities:
# Issues and Questions

# Communities: Issues and Questions

- **Some issues with community detection:**

  - Many different formalizations of clustering objective functions

  - Objectives are NP-hard to optimize exactly

  - Methods can find clusters that are systematically "biased"

- **Questions:**

  - **How well do algorithms optimize objectives?**

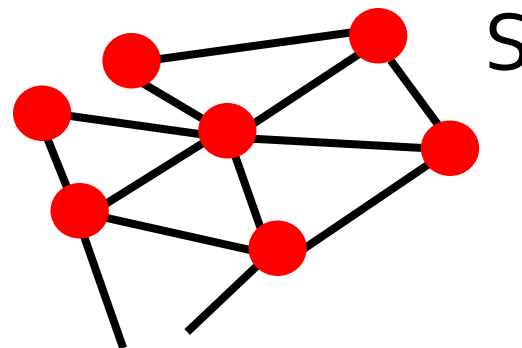  - **What clusters do different methods find?**

# Many Different Objective Functions

**Single-criterion:**

- Modularity: $m-E(m)$
- Edges cut: $c$

**Multi-criterion:**

- <u>Conductance</u>: $c/(2m+c)$
- Expansion: $c/n$
- Density: $1-m/n^2$
- CutRatio: $c/n(N-n)$
- Normalized Cut: $c/(2m+c) + c/2(M-m)+c$
- Flake-ODF: *frac. of nodes with more than ½ edges pointing outside S*

S

$n$: nodes in S
$m$: edges in S
$c$: edges pointing
outside S

# Many Classes of Algorithms

**Many algorithms to that implicitly or explicitly optimize objectives and extract communities:**
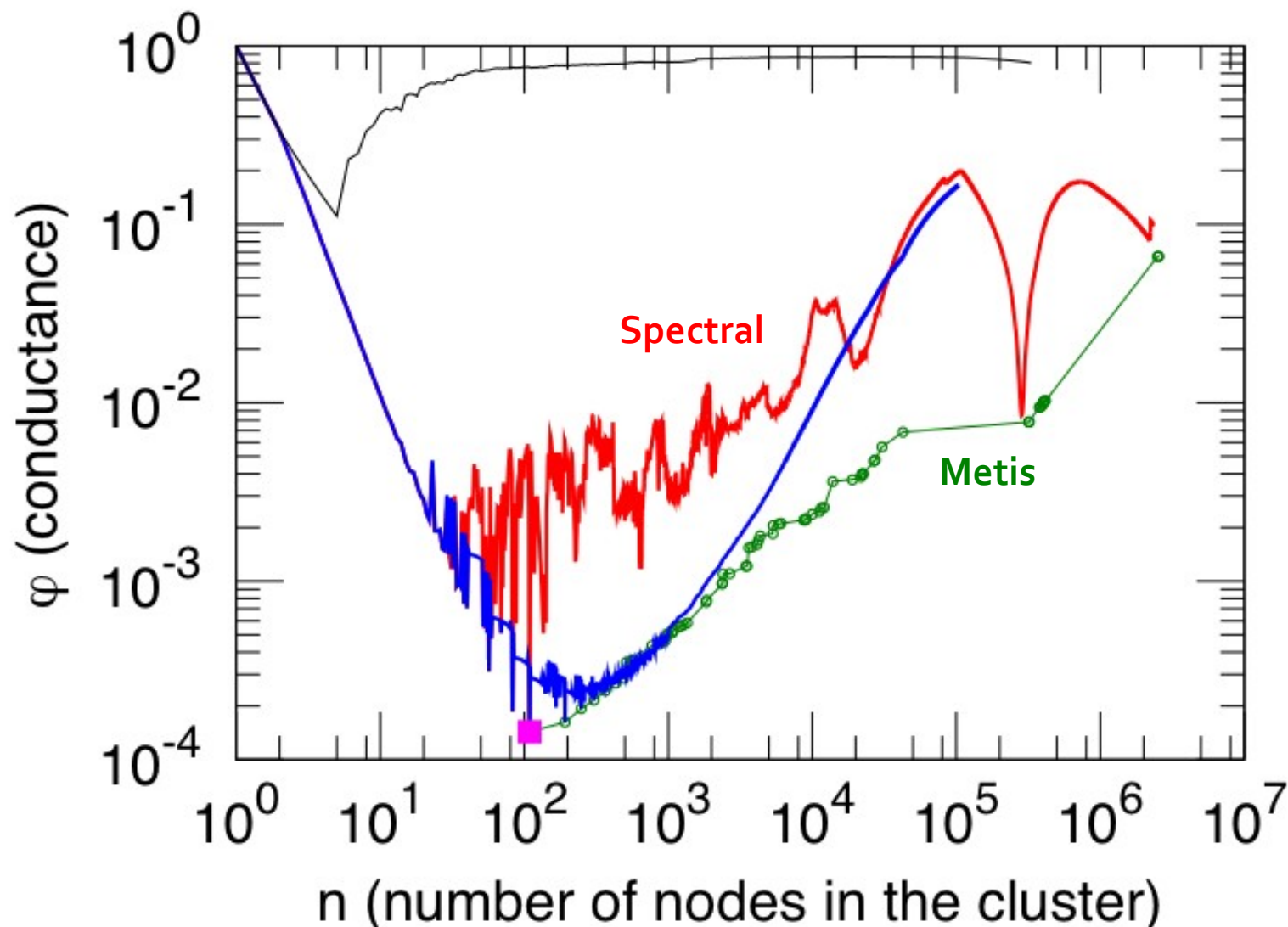
- **Heuristics:**
  - Girvan-Newman, Modularity optimization: popular heuristics
  - Metis: multi-resolution heuristic [Karypis-Kumar '98]
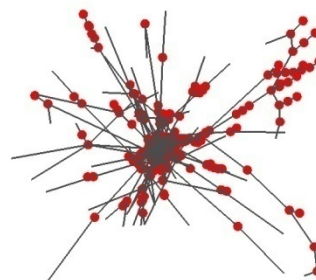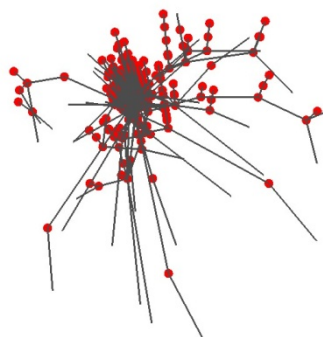
- **Theoretical approximation algorithms:**
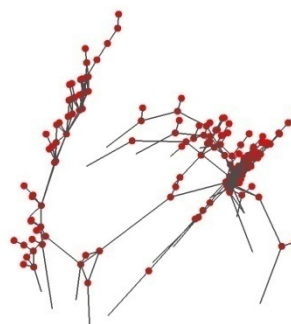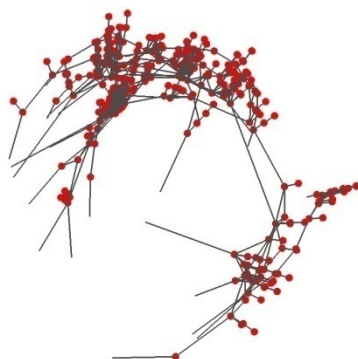  - Spectral partitioning

# NCP: Live Journal

# Properties of Clusters (1)

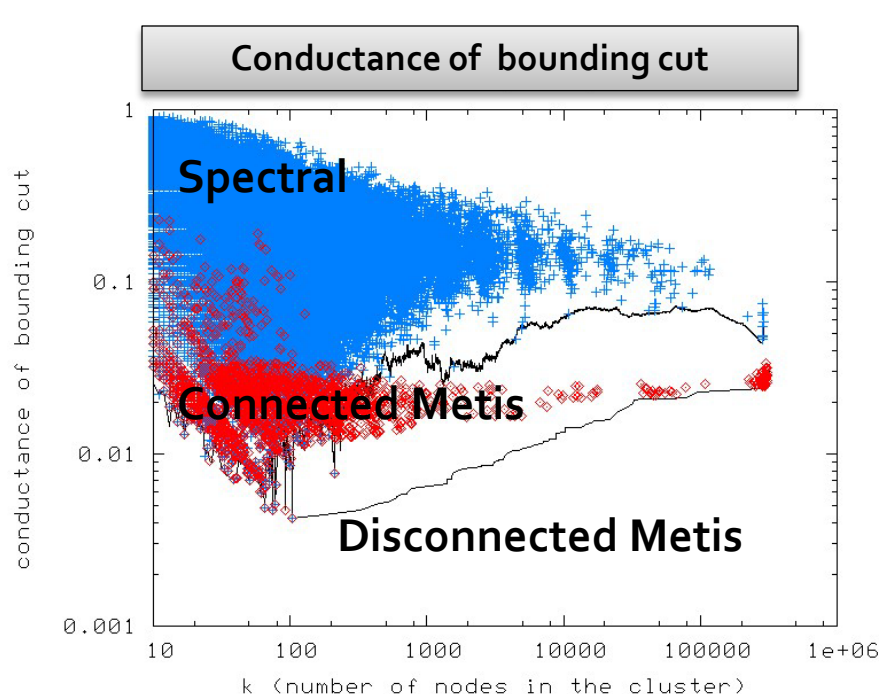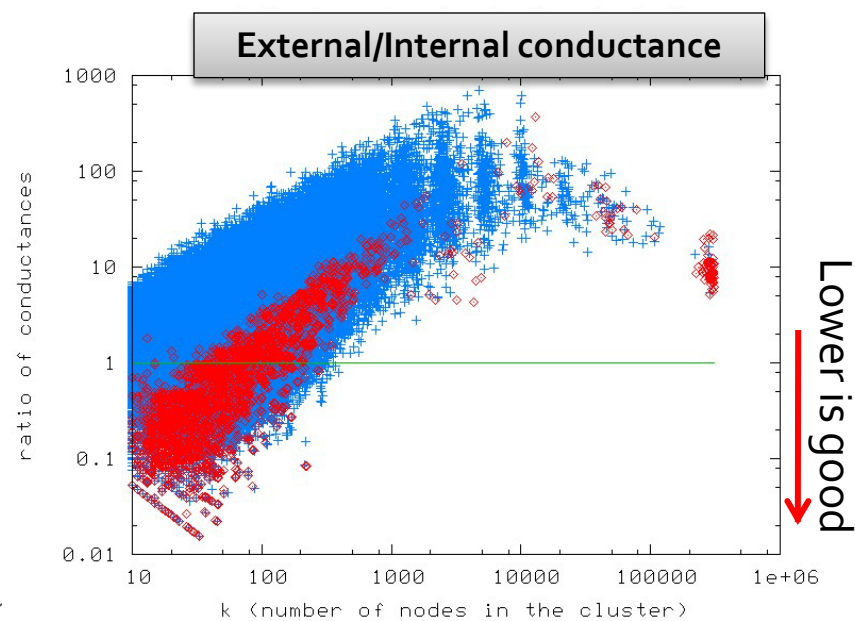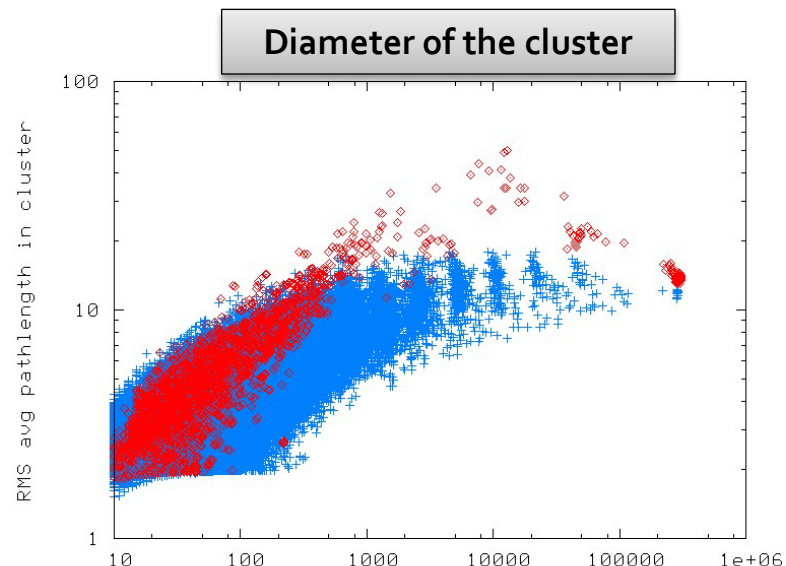## 500 node communities from Spectral:

## 500 node communities from Metis:

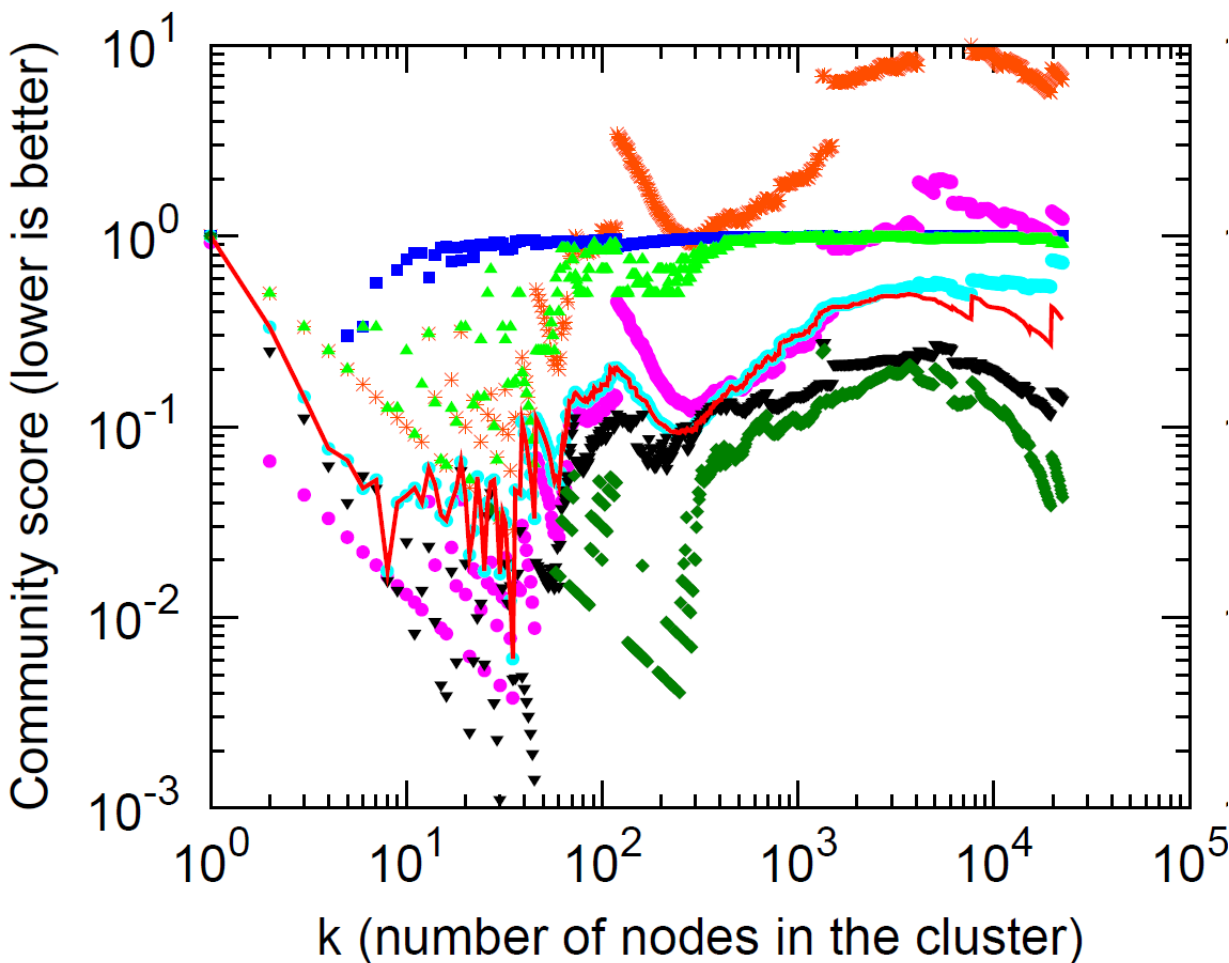# Properties of Clusters (2)

**Conductance of bounding cut**



Spectral

Connected Metis

Disconnected Metis

**Diameter of the cluster**



**External/Internal conductance**



Lower is good

- Metis (red) gives sets with better conductance

- Spectral (blue) gives tighter and more well-rounded sets

# Multi-criterion Objectives



- **All qualitatively similar**
- **Observations:**
  - Conductance, Expansion, Norm-cut, Cut-ratio are similar
  - Flake-ODF prefers larger clusters
  - Density is bad
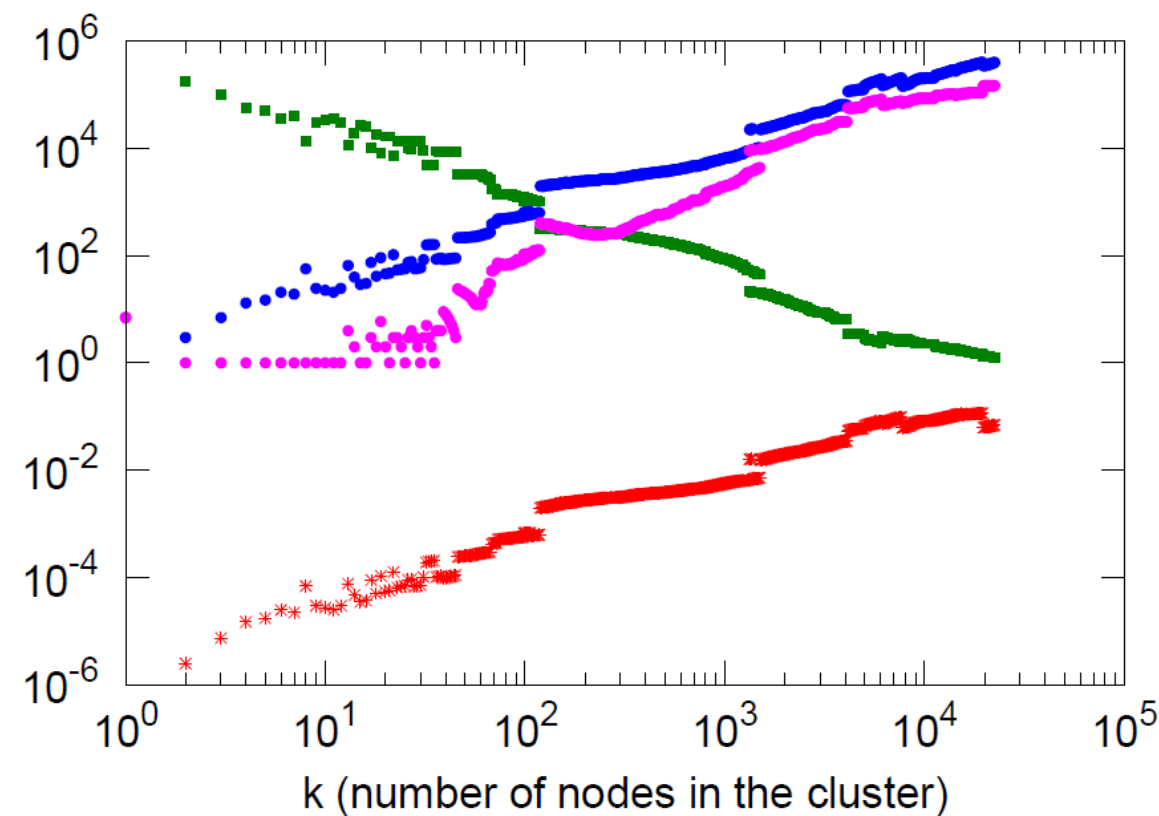  - Cut-ratio has high variance

Legend:
- Conductance — (red line)
- Expansion — * (asterisk)
- Internal Density — ■ (blue square)
- Cut Ratio — ● (magenta circle)
- Normalized Cut — ● (cyan circle)
- Maximum ODF — ▲ (green triangle)
- Avg ODF — ▼ (black triangle)
- Flake ODF — ◆ (green diamond)

# Single-criterion Objectives

**Observations:**

- All measures are monotonic
- **Modularity**
  - prefers large clusters
  - Ignores small clusters

k (number of nodes in the cluster)

Modularity  *  Modularity Ratio  ■  Volume  ●  Edges cut  ●